

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2011-221979

(P2011-221979A)

(43) 公開日 平成23年11月4日(2011.11.4)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 17/30 (2006.01)	G06F 17/30 350C	5B075
	G06F 17/30 210A	
	G06F 17/30 110C	

審査請求 未請求 請求項の数 5 O L (全 14 頁)

(21) 出願番号	特願2010-152556 (P2010-152556)	(71) 出願人	502192546 清華大学 中華人民共和国北京市海淀区清華大学 郵 編 100084
(22) 出願日	平成22年7月3日 (2010.7.3)	(71) 出願人	000155469 株式会社野村総合研究所 東京都千代田区丸の内一丁目6番5号
(31) 優先権主張番号	201010140447.0	(74) 代理人	100096002 弁理士 奥田 弘之
(32) 優先日	平成22年4月2日 (2010.4.2)	(74) 代理人	100091650 弁理士 奥田 規之
(33) 優先権主張国	中国 (CN)	(72) 発明者	李 春平 中華人民共和国北京市海淀区清華大学内
		(72) 発明者	王 益▲びん▼ 中華人民共和国北京市海淀区清華大学内 最終頁に続く

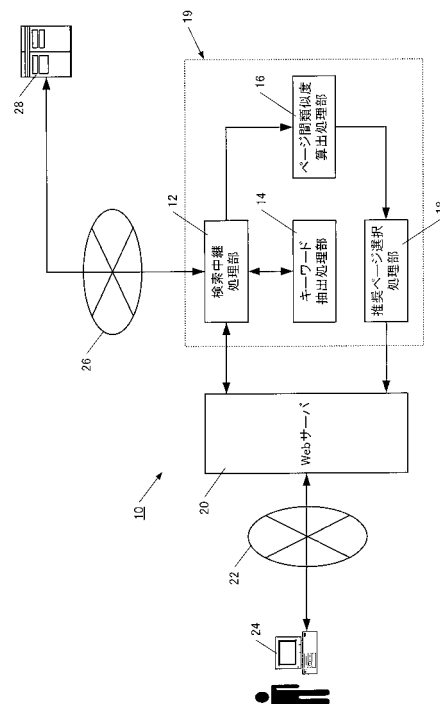
(54) 【発明の名称】 検索支援システム、検索支援方法及び検索支援プログラム

(57) 【要約】 (修正有)

【課題】 ユーザの検索意図に合致したWebページを効率的に提示する。

【解決手段】 検索キーワードを検索サーバ28に送信して検索を依頼する手段と、送信された検索結果リストをクライアント端末24に送信し、基準Webページの選択を促す手段と、この基準Webページ中のテキストを形態素単元に分解し、名詞をキーワードとして抽出する手段と、各抽出キーワードと検索キーワードを検索サーバ28に送信して検索を依頼する手段と、検索結果リストの中から、それぞれページランク順に上位20件のWebページを類似候補ページとして抽出する手段と、各検索式に係る類似候補ページ間における類似度を算出する手段と、類似度の高い方から20件の検索式に係る抽出キーワードを重要語と認定する手段と、全類似候補ページ中で重要語を3種以上含むものを類似ページと認定する手段と、類似ページリスト画面をクライアント端末24に送信する手段を備える。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

クライアント端末から送信された検索キーワードを検索サーバに送信し、検索を依頼する手段と、

上記検索サーバから送信された検索結果リストが記載された画面を上記クライアント端末に送信し、基準Webページの選択を促す手段と、

上記クライアント端末から基準Webページの選択情報を受信した場合に、この基準Webページ中のテキストを形態素単位に分解し、特定の品詞に係るキーワードを抽出する手段と、

これらの抽出キーワードと上記検索キーワードとをAND条件で繋いだ検索式を抽出キーワード毎に生成し、これらの検索式を上記検索サーバに送信して検索を依頼する手段と、

検索サーバから送信された各検索式に係る検索結果リストの中から、それぞれページランク順に上位所定件数のWebページを類似候補ページとして抽出する手段と、

各検索式に係る類似候補ページ間における類似度を算出するページ間類似度算出手段と、

類似度の高い所定件数の検索式に係る抽出キーワードを、重要語と認定する手段と、

上記の全類似候補ページの中で、上記重要語を所定種類以上含むものを類似ページと認定する手段と、

この類似ページのリストが記載された画面を生成し、上記クライアント端末に送信する手段と、

を備えたことを特徴とする検索支援システム。

【請求項 2】

上記ページ間類似度算出手段は、

各検索式に係る類似候補ページ中の一の類似候補ページを比較対象ページとして設定する処理と、この比較対象ページと残りの類似候補ページ間の類似度を個別に算出する処理と、各算出結果の中で類似度が上位所定数のものを抽出し、これらの類似度の平均値を隣接値として算出する処理を、当該検索式に係る全ての類似候補ページが比較対象ページとして設定されるまで繰り返した後、

得られた隣接値の中で最も大きな値の隣接値を当該検索式に係るページ間類似度と認定することを特徴とする請求項 1 に記載の検索支援システム。

【請求項 3】

上記ページ間類似度算出手段は、

上記の比較対象ページと、他の類似候補ページを形態素単位に分解し、各ページから所定の品詞に係る形態素を抽出する処理と、抽出された各形態素のTF-IDF値を算出する処理と、この各形態素のTF-IDF値に基づいて各ページをベクトル化する処理と、比較対象ページのベクトルと他の類似候補ページのベクトルとの内積値を、両ページ間の類似度として算出する処理を実行することを特徴とする請求項 2 に記載の検索支援システム。

【請求項 4】

クライアント端末から送信された検索キーワードを検索サーバに送信し、検索を依頼するステップと、

上記検索サーバから送信された検索結果リストが記載された画面を上記クライアント端末に送信し、基準Webページの選択を促すステップと、

上記クライアント端末から基準Webページの選択情報を受信した場合に、この基準Webページ中のテキストを形態素単位に分解し、特定の品詞に係るキーワードを抽出するステップと、

これらの抽出キーワードと上記検索キーワードとをAND条件で繋いだ検索式を抽出キーワード毎に生成し、これらの検索式を上記検索サーバに送信して検索を依頼するステップと、

検索サーバから送信された各検索式に係る検索結果リストの中から、それぞれページランク順に上位所定件数のWebページを類似候補ページとして抽出するステップと、

10

20

30

40

50

各検索式に係る類似候補ページ間における類似度を算出するページ間類似度算出ステップと、

類似度の高い所定件数の検索式に係る抽出キーワードを、重要語と認定するステップと、

上記の全類似候補ページの中で、上記重要語を所定種類以上含むものを類似ページと認定するステップと、

この類似ページのリストが記載された画面を生成し、上記クライアント端末に送信するステップと、

からなることを特徴とする検索支援方法。

【請求項 5】

10

コンピュータを、

クライアント端末から送信された検索キーワードを検索サーバに送信し、検索を依頼する手段、

上記検索サーバから送信された検索結果リストが記載された画面を上記クライアント端末に送信し、基準Webページの選択を促す手段、

上記クライアント端末から基準Webページの選択情報を受信した場合に、この基準Webページ中のテキストを形態素単位に分解し、特定の品詞に係るキーワードを抽出する手段、

これらの抽出キーワードと上記検索キーワードとをAND条件で繋いだ検索式を抽出キーワード毎に生成し、これらの検索式を上記検索サーバに送信して検索を依頼する手段、

検索サーバから送信された各検索式に係る検索結果リストの中から、それぞれページランク順に上位所定件数のWebページを類似候補ページとして抽出する手段、

20

各検索式に係る類似候補ページ間における類似度を算出するページ間類似度算出手段、

類似度の高い所定件数の検索式に係る抽出キーワードを、重要語と認定する手段、

上記の全類似候補ページの中で、上記重要語を所定種類以上含むものを類似ページと認定する手段、

この類似ページのリストが記載された画面を生成し、上記クライアント端末に送信する手段、

として機能させることを特徴とする検索支援プログラム。

【発明の詳細な説明】

【技術分野】

30

【0001】

この発明は検索支援システム、検索支援方法及び検索支援プログラムに係り、特に、ユーザの意図に合致したWebページを効率的に抽出する技術に関する。

【背景技術】

【0002】

今日、Google（登録商標）やYahoo!（登録商標）といったインターネット上の検索サイトに検索キーワードを送信することにより、誰でもが手軽に様々な情報を入手することができる。

例えば、ある銘柄の株式の購入を検討している一般投資家が、当該銘柄の企業名を検索サイトの検索窓に入力して検索をリクエストすると、当該企業に関するWebページのリストが検索結果として送信され、Webブラウザ上に表示される。

40

これに対し投資家は、ニュース記事やリリース記事、ブログ記事等を次々と閲覧し、当該企業の最近の動向をチェックする。

そして、新製品情報や不祥事情報など、株価に影響を及ぼすような記事を見出した場合には、当該記事中に用いられている適当なキーワードを企業名に追加し、さらに検索を続けることで、必要な情報を集中的に収集することが可能となる。

【非特許文献 1】Google インターネットURL:<http://www.google.co.jp/> 検索日：平成 22 年 1 月 16 日

【非特許文献 2】Yahoo! JAPAN インターネットURL:<http://www.yahoo.co.jp/> 検索日：平成 22 年 1 月 16 日

50

【発明の開示】

【発明が解決しようとする課題】

【0003】

しかしながら、このような従来の検索サイトを用いた検索方式の場合、ユーザには追加のキーワードを自分で選定して検索サイトに再投入し、イメージ通りの検索結果が得られるまで何度も操作を繰り返す手間がかかっていた。

この発明は、従来のこのような問題を解決するために案出されたものであり、ユーザの検索意図に合致したWebページを効率的に提示することができる技術の実現を目的としている。

【課題を解決するための手段】

【0004】

上記の目的を達成するため、請求項1に記載した検索支援システムは、クライアント端末から送信された検索キーワードを検索サーバに送信し、検索を依頼する手段と、上記検索サーバから送信された検索結果リストが記載された画面を上記クライアント端末に送信し、基準Webページの選択を促す手段と、上記クライアント端末から基準Webページの選択情報を受信した場合に、この基準Webページ中のテキストを形態素単位に分解し、特定の品詞に係るキーワードを抽出する手段と、これらの抽出キーワードと上記検索キーワードとをAND条件で繋いだ検索式を抽出キーワード毎に生成し、これらの検索式を上記検索サーバに送信して検索を依頼する手段と、検索サーバから送信された各検索式に係る検索結果リストの中から、それぞれページランク順に上位所定件数のWebページを類似候補ページとして抽出する手段と、各検索式に係る類似候補ページ間における類似度を算出するページ間類似度算出手段と、類似度の高い所定件数の検索式に係る抽出キーワードを、重要語と認定する手段と、上記の全類似候補ページの中で、上記重要語を所定種類以上含むものを類似ページと認定する手段と、この類似ページのリストが記載された画面を生成し、上記クライアント端末に送信する手段とを備えたことを特徴としている。

【0005】

請求項2に記載した検索支援システムは、請求項1のシステムを前提とし、さらに上記ページ間類似度算出手段が、各検索式に係る類似候補ページ中の一の類似候補ページを比較対象ページとして設定する処理と、この比較対象ページと残りの類似候補ページ間の類似度を個別に算出する処理と、各算出結果の中で類似度が上位所定数のものを抽出し、これらの類似度の平均値を隣接値として算出する処理を、当該検索式に係る全ての類似候補ページが比較対象ページとして設定されるまで繰り返した後、得られた隣接値の中で最も大きな値の隣接値を当該検索式に係るページ間類似度と認定することを特徴としている。

【0006】

請求項3に記載した検索支援システムは、請求項2のシステムを前提とし、さらに上記ページ間類似度算出手段が、上記の比較対象ページと、他の類似候補ページを形態素単位に分解し、各ページから所定の品詞に係る形態素を抽出する処理と、抽出された各形態素のTF-IDF値を算出する処理と、この各形態素のTF-IDF値に基づいて各ページをベクトル化する処理と、比較対象ページのベクトルと他の類似候補ページのベクトルとの内積値を、両ページ間の類似度として算出する処理を実行することを特徴としている。

【0007】

請求項4に記載した検索支援方法は、クライアント端末から送信された検索キーワードを検索サーバに送信し、検索を依頼するステップと、上記検索サーバから送信された検索結果リストが記載された画面を上記クライアント端末に送信し、基準Webページの選択を促すステップと、上記クライアント端末から基準Webページの選択情報を受信した場合に、この基準Webページ中のテキストを形態素単位に分解し、特定の品詞に係るキーワードを抽出するステップと、これらの抽出キーワードと上記検索キーワードとをAND条件で繋いだ検索式を抽出キーワード毎に生成し、これらの検索式を上記検索サーバに送信して検索を依頼するステップと、検索サーバから送信された各検索式に係る検索結果リストの中から、それぞれページランク順に上位所定件数のWebページを類似候補ページとして抽出

10

20

30

40

50

するステップと、各検索式に係る類似候補ページ間における類似度を算出するページ間類似度算出ステップと、類似度の高い所定件数の検索式に係る抽出キーワードを、重要語と認定するステップと、上記の全類似候補ページの中で、上記重要語を所定種類以上含むものを類似ページと認定するステップと、この類似ページのリストが記載された画面を生成し、上記クライアント端末に送信するステップとからなることを特徴としている。

【0008】

請求項5に記載した検索支援プログラムは、コンピュータを、クライアント端末から送信された検索キーワードを検索サーバに送信し、検索を依頼する手段、上記検索サーバから送信された検索結果リストが記載された画面を上記クライアント端末に送信し、基準Webページの選択を促す手段、上記クライアント端末から基準Webページの選択情報を受信した場合に、この基準Webページ中のテキストを形態素単位に分解し、特定の品詞に係るキーワードを抽出する手段、これらの抽出キーワードと上記検索キーワードとをAND条件で繋いだ検索式を抽出キーワード毎に生成し、これらの検索式を上記検索サーバに送信して検索を依頼する手段、検索サーバから送信された各検索式に係る検索結果リストの中から、それぞれページランク順に上位所定件数のWebページを類似候補ページとして抽出する手段、各検索式に係る類似候補ページ間における類似度を算出するページ間類似度算出手段、類似度の高い所定件数の検索式に係る抽出キーワードを、重要語と認定する手段、上記の全類似候補ページの中で、上記重要語を所定種類以上含むものを類似ページと認定する手段、この類似ページのリストが記載された画面を生成し、上記クライアント端末に送信する手段として機能させることを特徴としている。

10

20

【発明の効果】

【0009】

請求項1に記載の検索支援システム、請求項4に記載の検索支援方法及び請求項5に記載の検索支援プログラムによれば、ユーザは検索結果リスト中から自己の検索目的に最も近い一のWebページを基準Webページとして選択するだけで、当該Webページに類似したWebページのリストを取得することができるため、キーワードの選定及び再検索リクエストの煩わしさからユーザを解放することが可能となる。

【0010】

請求項2に記載の検索支援システムの場合、検索キーワードと各抽出キーワードとの組合せ単位のページ間類似度を算出するに際し、まず一対の類似候補ページ同士で類似度を次々と算出していき、その過程で極端に低い値の類似度を排除した上で平均値を求め、各平均値の中で最大のものをページ間類似度として採用する方式を採用しているため、内容の空疎なノイズ的なページの影響を廃して、より信頼性の高いページ間類似度を導くことが可能となる。

30

【0011】

請求項3に記載の検索支援システムの場合、両Webページに含まれる用語の構成や出現頻度に基づいて具体的な類似度が算出される仕組みであるため、記述内容に基づく類似性を算出結果に正確に反映させることが可能となる。

【発明を実施するための最良の形態】

【0012】

図1は、この発明に係る検索支援システム10の機能構成を示すブロック図であり、検索中継処理部12と、キーワード抽出処理部14と、ページ間類似度算出処理部16と、推奨ページ選択処理部18とを備えた検索支援サーバ19と、Webサーバ20とから構成される。

40

上記の検索中継処理部12、キーワード抽出処理部14、ページ間類似度算出処理部16、推奨ページ選択処理部18は、検索支援サーバ19のCPUが、OS及びアプリケーションプログラムに従って必要な処理を実行することによって実現される。

【0013】

検索中継処理部12及び推奨ページ選択処理部18は、Webサーバ20及びインターネット22を介して、ユーザの操作するクライアント端末24と接続される。

また、検索中継処理部12は、インターネット26を介して、Google（登録商標）やYahoo!

50

(登録商標)等の検索サーバ28と接続されている。この検索サーバ28には、インデックス情報を参照してユーザが入力した検索キーワードを含むWebページを抽出すると共に、抽出したWebページを所定のアルゴリズムに従ってランク付け(順位付け)する機能を備えた検索エンジンが搭載されている。

【0014】

つぎに、図2のフローチャートに従い、この検索支援システム10における全体的な処理手順を説明する。

まず、一般ユーザがクライアント端末24を操作してWebサーバ20内の検索支援サイトにアクセスし、Webブラウザ上に表示された検索窓に検索キーワードを入力して検索ボタンをクリックすると、クライアント端末24からWebサーバ20に検索キーワードが送信される。

10

【0015】

この検索キーワードをWebサーバ20経由で受信した検索中継処理部12は(S10)、これを検索サーバ28に送信し、検索を依頼する(S12)。

そして、検索サーバ28から検索結果リストが送信されると、これを受けた検索中継処理部12は(S14)、検索結果リストをWebサーバ20に送信する。

Webサーバ20は、この検索結果リストが記述された画面(Htmlファイル)を生成し、クライアント端末24に送信する(S16)。

【0016】

この結果、図3に示すように、検索結果リスト画面30がクライアント端末24のWebブラウザ上に表示される。

20

この図においては、「野村総合研究所」という検索キーワードの投入に対して、「野村総合研究所(NRI)」、「採用情報」、「NRI-Wikipedia」等のタイトルと、それぞれの概要情報が、検索結果リストとして列挙されている。

【0017】

これに対しユーザは、各タイトルをクリックして対応のWebページをWebブラウザ上に表示させ、個別に内容をチェックする。

そして、注目すべき内容を備えたWebページを発見した場合、上記の検索結果リスト画面30に戻り、当該Webページのタイトル近傍に設けられたチェックボックス32にチェックを入れた上で、「類似ページを見る」ボタン34をクリックする。

30

この結果、クライアント端末24からWebサーバ20に対して、基準Webページの選択情報を伴う類似ページの検索リクエストが送信される。

【0018】

Webサーバ20からこの基準Webページの選択情報を受け取った検索中継処理部12は(S18)、基準WebページのURLにアクセスして当該Webページを取得し(S19)、キーワード抽出処理部14に渡す。

これを受けたキーワード抽出処理部14は、当該基準ページに含まれるテキストに対して形態素解析処理を実行する(S20)。

ここで「形態素解析処理」とは、自然言語で記述された文を、意味を有する最小の言語単位である形態素に分解し、それぞれの品詞を特定する処理をいう。

40

【0019】

つぎにキーワード抽出処理部14は、各形態素の中から特定の品詞に係るものを抽出する(S22)。ここでは、名詞の形態素が90個抽出されたものとして、話を進める。

【0020】

つぎに、キーワード抽出処理部14から抽出キーワードを渡された検索中継処理部12は、ユーザが最初に入力した検索キーワードと、各抽出キーワードとをAND条件で繋いだ90個分の検索式を生成し、各検索式を検索サーバ28に送信して検索を依頼する(S24)。

図4は、この場合の具体例を示すものであり、当初の検索キーワードである「野村総合研究所」に対して、「本日」、「新サービス」、「発表」等の抽出キーワードがスペースを介してAND条件で結ばれた状態で、検索サーバ28に投入されるイメージを表している。

50

【 0 0 2 1 】

そして、検索サーバ28から「検索キーワード&抽出キーワード」単位での検索結果を受信すると（S26）、検索中継処理部12はページランクが上位20件以内のWebページを類似候補ページとして抽出し（S28）、ページ間類似度算出処理部16に渡す。

図4には、90組の検索式（検索キーワード&抽出キーワード）についてそれぞれ20件のWebページが抽出された結果、20（件）×90（組）=1,800（ページ）の類似候補ページ群が得られた例が示されている。

【 0 0 2 2 】

上記のページランクは、検索結果である各Webページの被リンク数やリンク元ページのページランク等に基づき、固有のアルゴリズムに従って検索サーバ28が付与するものであり、各Webページの有用性を示す指標と評価できる。

このため、このシステム10においては、ページランクの上位20件を類似候補ページとして抽出することにより、有用性の低いWebページを排除している。ただし、この抽出件数の「20件」は一例であり、他の閾値（件数）を適用することも当然に可能である。また、固定的な件数の代わりに、所定のページランク以上のWebページを全て類似候補ページとして抽出することもできる。

【 0 0 2 3 】

検索中継処理部12から類似候補ページ群を渡されたページ間類似度算出処理部16は、検索式（検索キーワード&抽出キーワード）単位で、上位20件の類似候補ページ間における類似度を算出する（S30）。このページ間類似度の具体的な算出手順については、後に詳述する。

【 0 0 2 4 】

ページ間類似度算出処理部16から、合計90組分のページ間類似度の算出結果を受け取った推奨ページ選択処理部18は、各ページ間類似度の数値を高い順に整列させ、上位20組の検索式に係る抽出キーワードを「重要語」と認定する（S32）。図4においては、90組の検索式から、「新サービス」、「アメリカ企業」、「コスト削減」等が重要語として抽出された例が示されている。

なお、この「上位20組」も一例であり、他の件数の抽出キーワードを重要語として抽出してもよい。

【 0 0 2 5 】

つぎに推奨ページ選択処理部18は、1,800ページに及ぶ類似候補ページ群の中から、上記20個の重要語の中で3種以上の重要語を含むものを類似ページとして選定し、類似ページリストを生成する（S34）。

各重要語は、当初の検索キーワード（野村総合研究所）との組合せにおいて、基準Webページの内容を特徴付ける文字通り重要な意義を有する語であると考えられるため、より多くの種類の重要語を含むWebページは、それだけ基準Webページに対する類似度が高いものと評価できる。

ただし、「3種以上」の閾値は一例であり、他の値を閾値として設定することができる。また、複数種類の重要語を含むことが必要とされるため、一種類の重要語（例えば「コスト削減」）のみが100回登場するようなWebページは、類似ページと認定されることはない。

【 0 0 2 6 】

推奨ページ選択処理部18からこの類似ページリストを受け取ったWebサーバ20は、この類似ページリストを含む画面（Htmlファイル）を生成し、クライアント端末24に送信する（S36）。

【 0 0 2 7 】

この結果、図5に示すように、クライアント端末24のWebブラウザ上に類似ページリスト画面40が表示される。

この類似ページリスト中のタイトルをユーザがクリックすると、ユーザが最初に選択した基準Webページの記述内容と類似した内容を備えたWebページが表示されることとなる。

10

20

30

40

50

例えば、発明者が行った実証試験の結果では、適合率（基準Webページと関係のあるページが含まれる割合）が平均で86%、再現率（重要なページが含まれる割合）が平均で60%となり、かなり高い精度であることが確認された。

【0028】

上記においては、一のWebページが基準Webページとして選択された例を示したが、ユーザは、図3の検索結果リスト画面30において2以上のチェックボックス32にチェックを入れることにより、複数の基準Webページを選択することもできる。

この場合、キーワード抽出処理部14は各基準Webページを形態素に分解し、特定品詞（例えば名詞）に係るキーワードを各基準Webページから抽出した上で、これらを一括の抽出キーワード群となす。

後は、検索中継処理部12等により、図2のS24以下の処理が順次実行されることにより、類似ページリスト画面40が生成され、クライアント端末24に送信される。

【0029】

つぎに、図6のフローチャート及び図7の説明図に従い、ページ間類似度算出処理部16によるページ間類似度の算出手順について説明する。

まずページ間類似度算出処理部16は、20件の類似候補ページの中から一のページを、比較対象ページとして設定する（S30-01）。図7(a)では、ページAが最初の比較対象ページに設定された例が示されている。

【0030】

つぎにページ間類似度算出処理部16は、TF-IDF及びベクトル空間法を用いて、残り19件のページ（ページB～ページT）とページAとの間の類似度を個別に算出する（S30-02）。このTF-IDF及びベクトル空間法を用いた類似度の具体的な算出手順については、後に詳述する。

【0031】

つぎにページ間類似度算出処理部16は、算出された計19個の類似度を高い順に整列させ、上位15個の類似度の平均値を隣接値として算出する（S30-03）。図7(a)においては、ページAを比較対象とした場合の隣接値として「0.21」の数値が導かれている。

【0032】

つぎにページ間類似度算出処理部16は、図7(b)に示すように、次のページBを比較対象に設定した上で（S30-04/N、S30-01）、上記と同様の手順に従い、隣接値を求める（S30-02、S30-03）。

そして、図7(c)に示すように、最後のページTを比較対象に設定した上で、上記と同様の手順に従って隣接値を求め終えた時点で（S30-04/Y）、ページ間類似度算出処理部16は全20件の隣接値の中で最大の値を備えたものを、当該検索式に係るページ間類似度と認定する（S30-05）。

図7においては、ある検索式に係るページ間類似度として、「0.32」が導かれた例が示されている。

【0033】

つぎに、図8のフローチャート及び図9、図10の説明図に従い、TF-IDF及びベクトル空間法を用いた類似度の具体的な算出手順について説明する。

まずページ間類似度算出処理部16は、各類似候補ページに対して形態素解析を行い、特定品詞（例えば名詞）の用語のみを抽出する（S30-02-01）。

図9の例では、ページの「今日が締め切りだ。今日も徹夜かな。」から「今日/締め切り/今日/徹夜」の用語が、ページの「今日も煮干しだ。飽き飽きだ。」から「今日/煮干し」の用語が、ページの「今日は天気がよい。野球をしよう。」から「今日/天気/野球」の用語が、ページの「天気がよい。サッカーをしよう。」から「天気/サッカー」の用語がそれぞれ取り出されている。なお、図9及び図10においては、説明を単純化するためページ～ページの4件の類似候補ページに基づいてページ間類似度を算出する例が示されているが、類似候補ページの数が増えた場合も基本的な考え方は同じである。

10

20

30

40

50

【 0 0 3 4 】

つぎにページ間類似度算出処理部16は、各ページにおける各用語の頻度 (TF / Term Frequency) を算出する (S 30-02-02)。例えば、ページ における「今日」の頻度は「 2 」となる。

【 0 0 3 5 】

つぎにページ間類似度算出処理部16は、用語毎に当該用語を含むページ数 (DF / Document Frequency) を算出し (S 30-02-03)、DF辞書50に格納する (S 30-02-04)。例えば、ページ ~ において「今日」を含むページは 3 件あるため、「今日」のDFは「 3 」となる。

【 0 0 3 6 】

つぎにページ間類似度算出処理部16は、このDF辞書50に基づいて各ページをベクトル化する。

例えば、ページ の場合はDF辞書50に収録された用語の中、「今日」「締め切り」「徹夜」の 3 種類の用語を含んでいるため、ページ間類似度算出処理部16はこれらの用語のDFに基づいて、IDF (Inverse Document Frequency) 及びTF-IDFを求める。

【 0 0 3 7 】

まずページ間類似度算出処理部16は、以下のようにして各用語のIDFを算出する (S 30-02-05)。

$$\begin{aligned} \text{IDF (今日)} &= \log (\text{文書数} / \text{DF}) \\ &= \log (4 / 3) \end{aligned}$$

【 0 0 3 8 】

つぎにページ間類似度算出処理部16は、以下のようにして各用語のTF-IDFを算出する (S 30-02-06)。

$$\begin{aligned} \text{TF-IDF (今日)} &= \text{TF (今日)} \times \text{IDF (今日)} \\ &= 2 \times \log (4 / 3) = 0.25 \end{aligned}$$

同様の処理により、ページ間類似度算出処理部16は「締め切り」のTF-IDF = 0.60、「徹夜」のTF-IDF = 0.60を算出する。

【 0 0 3 9 】

ここで、ページ に含まれる「今日」「締め切り」「徹夜」の 3 種類の用語はDF辞書50における掲載順が 1 ~ 3 番であるため、図 1 0 に示すように、ベクトル要素として 1 ~ 3 行までに0.25、0.60、0.60の数値が代入され、他の用語の掲載順に対応する行には0.00が代入されたベクトルがページ間類似度算出処理部16によって生成され、ページ のベクトルとなされる (S 30-02-07)。

【 0 0 4 0 】

このページ のベクトル長は0.89であるため、ページ間類似度算出処理部16はベクトル長を 1 に揃えるための正規化処理を各数値に対して施し (S 30-02-08)、最終的に0.28、0.68、0.68、0.00、0.00、0.00、0.00の数値が充填されたベクトルを導く。

【 0 0 4 1 】

ページ の場合にはDF辞書50に収録された用語の中、「今日」「煮干し」の 2 種類の用語を含んでおり、これらの用語のDF辞書50における掲載順が 1 番と 4 番であるため、ベクトル要素として 1 行目及び 4 行目に0.12及び0.60の数値が代入され、他の用語の掲載順に対応する行には0.00が代入されたベクトルが生成された後 (S 30-02-07)、上記と同様の正規化処理により (S 30-02-08)、最終的に0.20、0.00、0.00、0.98、0.00、0.00、0.00の数値が充填されたベクトルが得られる。

【 0 0 4 2 】

また、文書 の場合はDF辞書50に収録された用語の中、「今日」「天気」「野球」の 3 種類の用語を含んでおり、これらの用語のDF辞書50における掲載順が 1 番と 5 番、6 番であるため、ベクトル要素として 1 行目、5 行目、6 行目にそれぞれ0.12、0.60、0.30の数値が代入され、他の用語の掲載順に対応する行には0.00が代入されたベクトルが生成された後 (S 30-02-07)、上記と同様の正規化処理により (S 30-02-08)、最終的に0.18、0.

10

20

30

40

50

00、0.00、0.00、0.88、0.44、0.00の数値が充填されたベクトルが得られる。

【0043】

また、ページ の場合はDF辞書50に収録された用語の中、「天気」「サッカー」の2種類の用語を含んでおり、これらの用語のDF辞書50における掲載順が6番と7番であるため、ベクトル要素として6行目及び7行目にそれぞれ0.30、0.60の数値が代入され、他の用語の掲載順に対応する行には0.00が代入されたベクトルが生成された後(S30-02-07)、上記と同様の正規化処理により(S30-02-08)、最終的に0.00、0.00、0.00、0.00、0.00、0.45、0.89の数値が充填されたベクトルが得られる。

【0044】

つぎにページ間類似度算出処理部16は、ページ のベクトルとページ のベクトルとの間の内積(距離)を求める(S30-02-09)。この内積値が、両ページ間の類似度となる。

以後、ページ間類似度算出処理部16は同様の手順に従い、ページ - ページ 間の類似度及びページ - ページ 間の類似度を算出する。

【図面の簡単な説明】

【0045】

【図1】この発明に係る検索支援システムの機能構成を示すブロック図である。

【図2】この検索支援システムの全体的な処理手順を示すフローチャートである。

【図3】検索結果リスト画面を示す図である。

【図4】類似ページの抽出に係る手順を示す概念図である。

【図5】類似ページリスト画面を示す図である。

【図6】ページ間類似度の算出手順を示すフローチャートである。

【図7】ページ間類似度の算出手順を示す概念図である。

【図8】TF-IDF及びベクトル空間法を用いた類似度の具体的な算出手順を示すフローチャートである。

【図9】TF-IDF及びベクトル空間法を用いた類似度の具体的な算出手順を示す説明図である。

【図10】TF-IDF及びベクトル空間法を用いた類似度の具体的な算出手順を示す説明図である。

【符号の説明】

【0046】

- 10 検索支援システム
- 12 検索中継処理部
- 14 キーワード抽出処理部
- 16 ページ間類似度算出処理部
- 18 推奨ページ選択処理部
- 19 検索支援サーバ
- 20 Webサーバ
- 22 インターネット
- 24 クライアント端末
- 26 インターネット
- 28 検索サーバ
- 30 検索結果リスト画面
- 32 チェックボックス
- 34 「類似ページを見る」ボタン
- 40 類似ページリスト画面
- 50 DF辞書

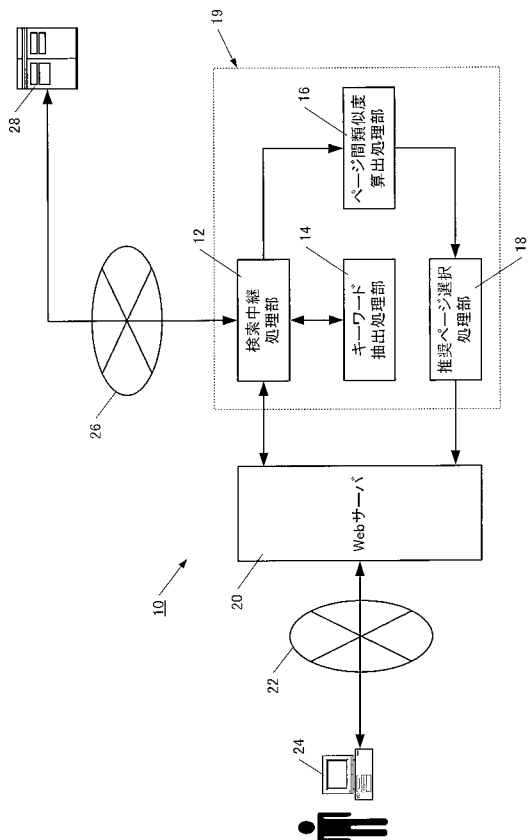
10

20

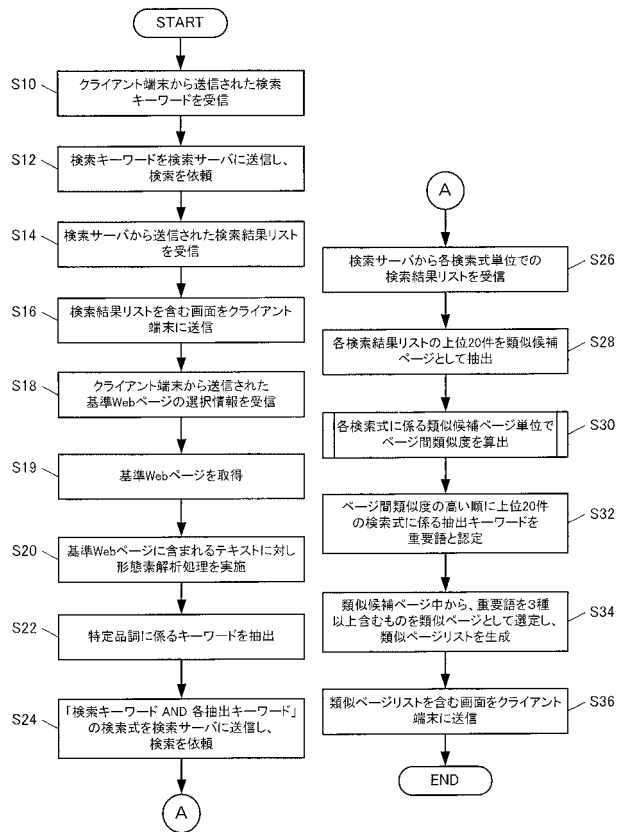
30

40

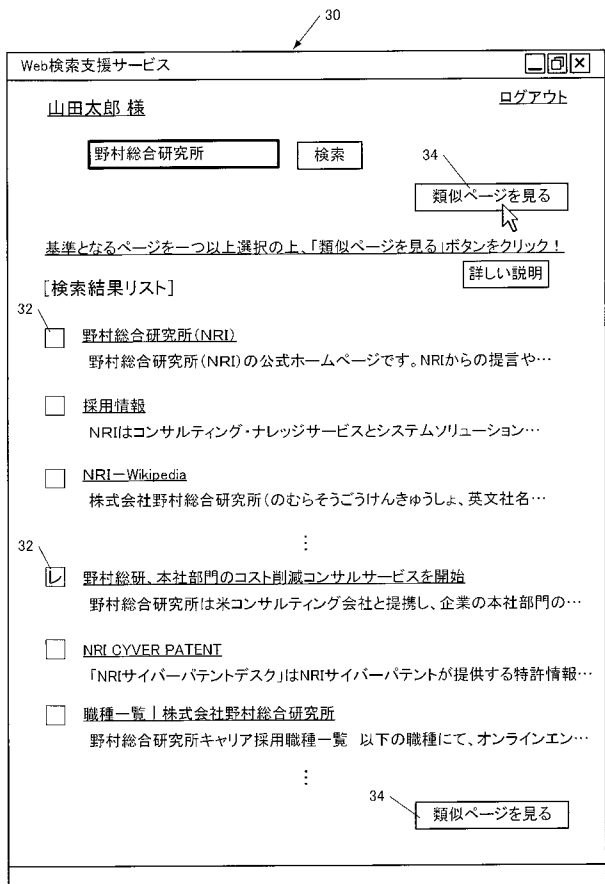
【図1】



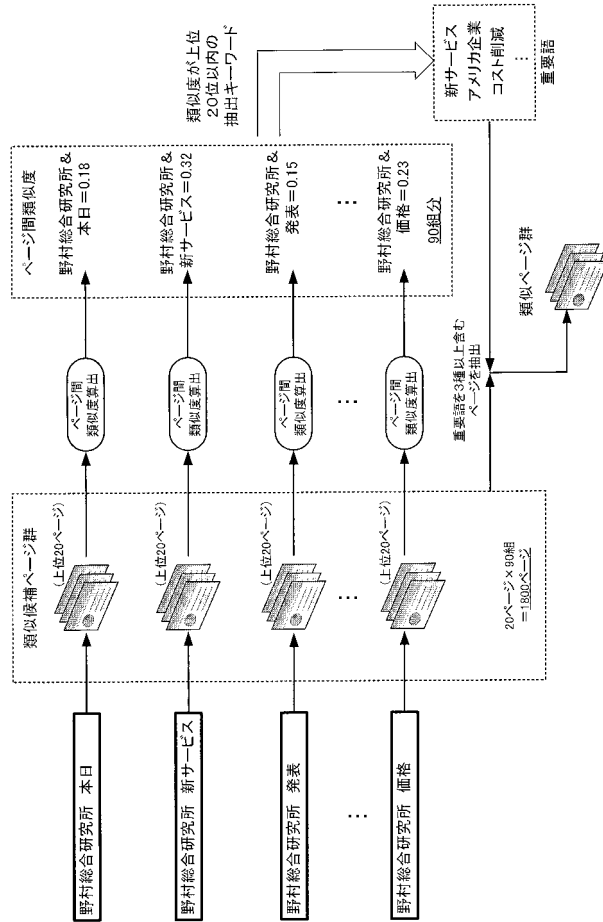
【図2】



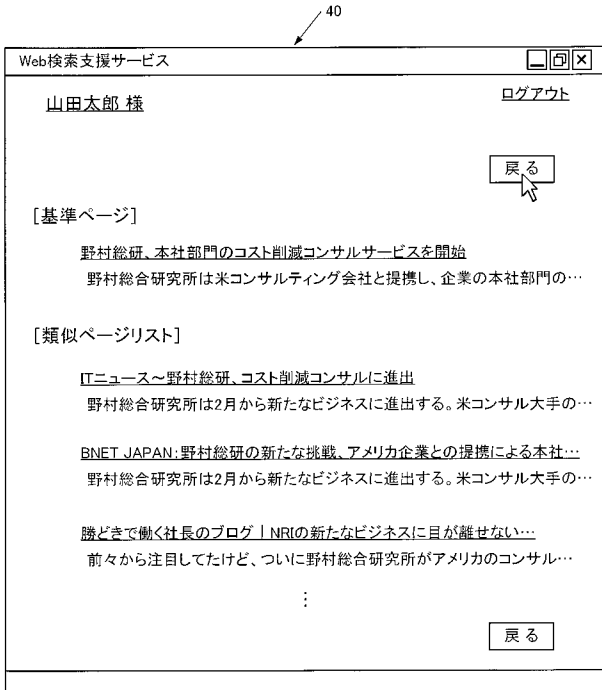
【図3】



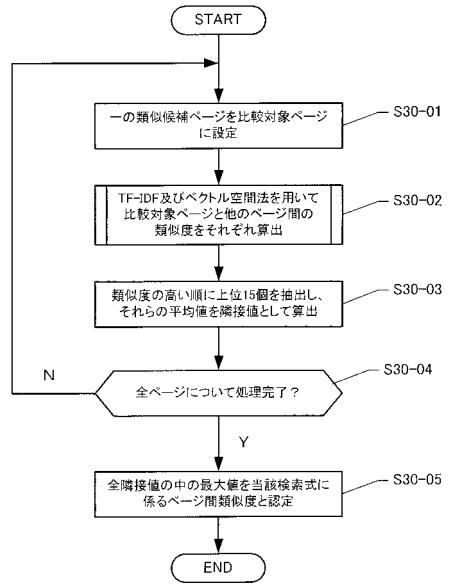
【図4】



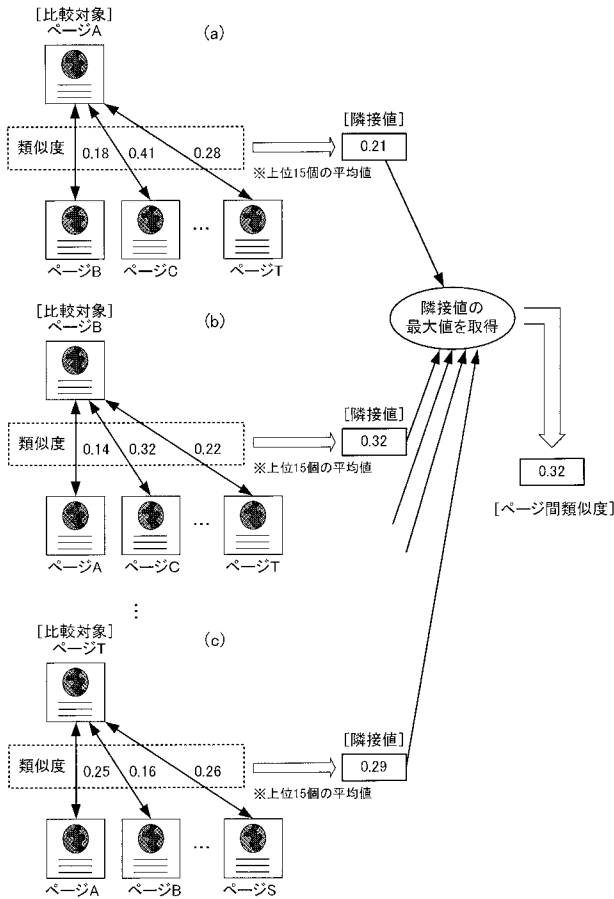
【 図 5 】



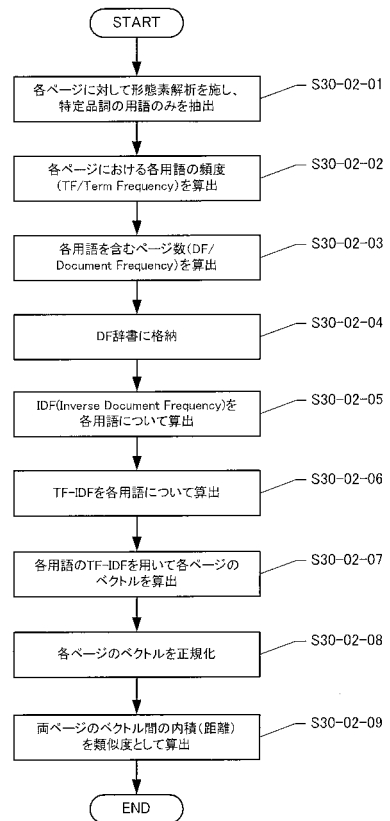
【 図 6 】



【 図 7 】



【 図 8 】



フロントページの続き

(72)発明者 阿部 昌平

東京都千代田区丸の内一丁目6番5号 株式会社野村総合研究所内

Fターム(参考) 5B075 KK02 ND03 NK32 QM05